

Joint Event Detection and Entity Resolution: a Virtuous Cycle

Matthias Gallé
Xerox Research Centre
Europe
6 Chemin de Maupertuis
38240 Meylan, France
matthias.galle@xrce.xerox.com

Jean-Michel Renders
Xerox Research Centre
Europe
6 Chemin de Maupertuis
38240 Meylan, France
jean-michel.renders@xrce.xerox.com

Guillaume Jacquet
Xerox Research Centre
Europe
6 Chemin de Maupertuis
38240 Meylan, France
guillaume.jacquet@xrce.xerox.com

ABSTRACT

Clustering web documents has numerous applications, such as aggregating news articles into meaningful events, detecting trends and hot topics on the Web, preserving diversity in search results, etc. At the same time, the importance of named entities and, in particular, the ability to recognize them and to solve the associated co-reference resolution problem are widely recognised as key enabling factors when mining, aggregating and comparing content on the Web.

Instead of considering these two problems separately, we propose in this paper a method that tackles jointly the problem of clustering news articles into events and cross-document co-reference resolution of named entities. The co-occurrence of named entities in the same clusters is used as an additional signal to decide whether two referents should be merged into one entity. These refined entities can in turn be used as enhanced features to re-cluster the documents and then be refined again, entering into a virtuous cycle that improves simultaneously the performances of both tasks. We implemented a prototype system and report results using the TDT5 collection of news articles, demonstrating the potential of our approach.

Categories and Subject Descriptors

H.2.8 [Database Applications]: Data Mining; H.3.3 [Information Search and Retrieval]: Clustering; I.2.7 [Natural Language Processing]: Text analysis; I.7 [Document and Text Processing]: Document Preparation

General Terms

Event-based information systems, natural-language processing, clustering

Keywords

Co-reference resolution, entity linking, topic detection and tracking, news aggregator

1. INTRODUCTION

Behind real-life web-oriented data mining applications are in most cases a suite of different modules, patched together in a sensible way. However, this ecosystem of algorithms is seldom acknowledged in research papers, which mostly focus on solving one single task, ignoring the general picture. On the one hand, this combination of algorithms can of course affect negatively the performance of the whole system, as errors tend to accumulate and serve as catalyzers of other errors. But on the other hand, this could also be used to boost the effectiveness of a single task, by resolving it together with another. Let us as an example a news article aggregator, the application that originated this research. Modern news aggregators combine multiple Natural Language Processing (NLP) components such as news article stemming/lemmatizing and parsing, Named Entity Recognition (NER), co-reference resolution, document clustering into news events¹, sentiment analysis, *etc.*. These components are executed independently, or in a pipeline where one component waits for the output from the other(s) in (for example) the order we listed above. Using a NER component and a co-reference resolution component for the clustering is straightforward as the named entities can be used as additional useful features when computing similarity measure between documents. In this paper, we analyze how the document clustering itself can improve the resolution of the cross-document named entity co-reference, which in turn is used to improve the clustering performance. This process can then be repeated, entering into a virtuous cycle that improves the resolution of both tasks.

Cross-document named entity co-reference consists in identifying in a set of documents all the textual expressions that refer to the same real entity referent (henceforth "entity" in this paper). It can be split into two sub-tasks, each with its own set of challenges. Intra-document co-reference resolution is more concerned with pronouns resolution, implicit mentions and linking correctly different noun phrases. Inter-document co-reference resolution on its side considers a set of documents and has to harmonize different or incomplete textual expressions referring to the same entity. It has to

¹In this paper, by "event" we mean "set of documents (or piece of documents) that refer to the same precise topic"

resolve spelling variants, use of parts of the complete name, acronyms, etc.

Cross-document named entity co-reference is particularly crucial in Information Extraction systems and namely in news event tracking, where new named entities/entity referents appear all the time and where there is a special interest in relating events through the actors (people and organizations) that took part in them and through the locations they occur in. Such news event tracking systems however have an advantage that has not been yet exploited in the literature. It creates clusters of news articles talking about specific events, providing social², geographical and temporal co-occurrences of the entities mentioned in the articles. This makes it possible to relate different textual expressions used by different sources for the same entity, which would otherwise slip through a standard named entity co-reference system.

We propose the use of co-occurrences in events as additional features (to be used in combination with other more traditional features, such as string similarities on the canonical surface form) in order to determine if two named entities co-refer to the same underlying entity. We use this building block in a system where this enhanced co-reference information can then be fed back to the clustering module in order to improve its quality. A general outline of such a system is depicted in Fig. 1. It receives as input a stream of documents. These are first preprocessed (html to text cleaning, lemmatization, *etc.*) and a NER system is run on each document individually. This NER system includes intra-document co-reference resolution, therefore the output is, for each document, a list of entities – represented by a *canonical string* – together with the number of times this entity appears (whatever the textual expression used) inside the document. The module is generally able to determine a canonical string, which is non-ambiguous, i.e. the most specific ones with respect to all mentions of an entity in a document. Such an intra-document co-reference component misses of course any possible variants due to spelling errors, transcription differences, *etc.*, besides those due to parser’s errors (such as taking a form (ADJ) (NOUN) as a named entity, e.g. "Faraway Andalusia"). Experimentally, we observed that state-of-the-art intra-document co-reference resolution algorithms very often tend to prioritize precision over recall, and over-generate entities. The phenomenon of co-reference is therefore much more common than the one of entity ambiguity, and our cross-document co-reference module addresses this by deciding which entities should be merged. In other words, we are considering the named entity co-reference problem as an alias detection task: we suppose that a unique real entity has several aliases (in particular because different sources of information use different aliases) and that, by mining the socio-temporal behavior of these aliases, we are able to merge two (or more) aliases as unique entity.

The first contribution of this paper is the use of social, geographical and temporal information in order to merge the named entities that co-refer. The second contribution is the

²in this paper, by "social co-occurrences" we mean co-occurrence of entities in the same event, i.e. in the same cluster of news articles

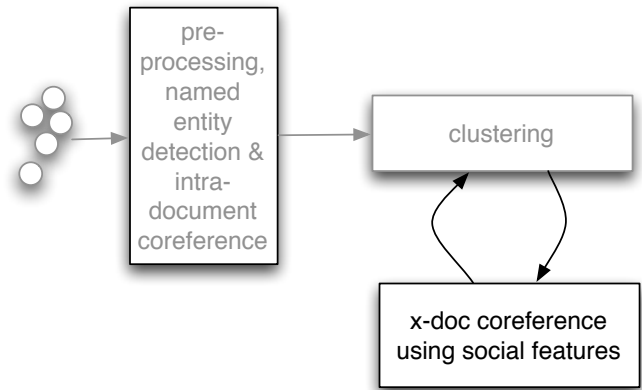


Figure 1: An overview of the system we propose. The documents arrive (in a stream, mini-batches or full-batch) and are preprocessed; named entities are extracted and intra-document co-reference resolution is applied. For the initial cross-document co-reference resolution, a trivial approach is used (two entities co-refer if there is an exact match between their canonical strings). Finally the documents are clustered and the information of co-occurrences in events is then used for an enhanced cross-document co-referencing. This is then fed back to the clustering algorithm and several iterations can be performed (fixed number or until convergence). Our novelty lies in the bold parts (the bootstrap cycle and the co-reference box).

creation of a virtuous cycle between the clustering and the cross-document co-reference module in order to iteratively improve the quality of both tasks.

2. RELATED WORK

This work stands in the intersection of two research fields in data mining, namely Topic Detection and Tracking (TDT) and cross-document co-reference resolution. While we did not find any previous work on tackling these two problems jointly, there is a rich literature considering these problems separately.

TDT [2] is concerned with monitoring news providers in order to extract *events*, which group news articles reporting the same concrete and precise topic. The most basic algorithm here is an incremental clustering algorithm that compares newly arrived articles to existing events and decides to assign it to one of the existing ones (Topic Tracking) or to create a new one (New Topic Detection). k -based clustering algorithms, where the number of clusters are predefined are unsuitable for this tasks due to the dynamic, constantly evolving nature of the newsphere. Several advanced clustering techniques have been proposed for this [16, 5] which instead take a similarity threshold as parameter, which controls how tight or sparse the cluster should be (corresponding to different levels of granularity of the events). In this work, we use the fully incremental (or single-pass) algorithm [17] that, while not the most effective algorithm for this task [5], runs fast and allows us to measure the impact of different parameters. This algorithm processes the

document one by one in chronological order. The first document creates a singleton cluster. A subsequent document is compared to the centroids of the existing clusters and the highest such similarity is recorded. If this similarity is above a pre-fixed threshold τ then the document gets assigned to this cluster and its centroids is updated. If not, it creates a singleton cluster by its own.

Classical works using statistical methods for named entity co-reference resolution traditionally compute a similarity between two mentions of an entity in the text, and then join some of these mentions into an equivalence cluster. This can be achieved for example with some kind of clustering: co-reference chains [9], tree traversal [8] or graph-cut algorithms [10]. The similarity between two mentions is computed using different features. A recent overview [15] for example, mentions 63 different features, and divides them into *lexical* (string-based comparison), *proximity* (number of words or paragraphs between two mentions), *grammatical* (based on POS, parse trees, etc) and *semantic* (gender, animacy, etc). Most of these features are based on the *content* of the named entity (e.g. some additional information provided by the parser on the superficial form it takes or on the local context given by surroundings words), while the *content-independent* features refer only to the close context and are always at the intra-document level.

A more global context has been used as feature for cross-document co-reference resolution. [6] uses the categories of the Open Directory Project and combines these with local context features (words in the same sentence). Such an approach is static by nature and no temporal information is used, which is a key feature in the domain of event detection.

A recent opposite approach [7] uses purely temporal information to solve the co-reference resolution problem: if two named entities have the same temporal profile (in other words, when their “bursts” of appearances are similar) and if their approximate string matching similarity is above some threshold, these named entities are claimed to refer to the same entity. Here, in addition to temporal information, we propose to introduce social and geographical contextual features. Those features are given by the clustering of news articles where each cluster of news articles is talking about one topic-coherent event, providing social, geographical and temporal co-occurrences of the entities mentioned by the articles.

Targeting the specific domain of the news-sphere, the European Media Monitor of the Joint Research Centre (JRC) of the EC has tackled the problem of multi-lingual named entity disambiguation. The main issue here is the correct translation of person names and the subsequent comparison between these translations. The similarity between two named entities is computed through the following process: transliteration into roman script, lowercasing, name normalization, vowel removal and finally edit distance [14]. An additional constraint mentioned in [12] is that, in order to merge two named entities, they have to appear in at least one identical news cluster. No other use of the social relationship is exploited. The same team released recently a database [14] of different name variants for the same entity.

This has been obtained with the process described above, plus using an external co-reference resource (Wikipedia) and manual revision. While this database is undoubtedly a valuable resource, it is tight to the news-sphere domain, has high-precision, low-recall rates and addresses only the most frequent named entities.

Our work is partially inspired by the approach proposed by Bhattacharya and Getoor [3, 4]. They target the use of complex entities in databases that provide additional information (like citation databases, where co-authors are linked, and different proceedings spelling can be linked through the use of the papers published therein). Their idea of using relational references to the same entity as additional information for merging shares some similarity with our use of temporal and context co-occurrences of different mentions. Note however the difference that our entities are much simpler (basically, just a canonical string) and the strong context information that can be provided through news events.

Another specificity of our proposed approach is to iteratively process both cross-document co-reference and event detection. A comparable iterative approach [3, 4] proposes a relational clustering algorithm for iteratively discovering entities by clustering references taking into account the clusters of co-occurring references. In this approach, clusters aggregate references (i.e. lexical expressions) while we are clustering documents. Their iterative approach aims at improving entity resolution task while we aim at improving both entity resolution (through cross-document co-reference task) and event detection tasks.

The opposed view that personal names are normally preserved between news articles of different languages and sources can be used for named entities discovery [13]. A related area where a similar idea has shown to be successful is *person name disambiguation*. The problem here is to disambiguate between two persons sharing the same name. The standard approach here [11] is to cluster the close context of each occurrence of the name.

3. JOINT EVENT RECOGNITION AND ENTITY RESOLUTION

Recall that we are jointly solving a document clustering problem – or more precisely, an event recognition and tracking task – and a cross-document co-reference problem. We adopt the extra assumption that an intra-document co-reference resolution step has been applied to the collection but, as it is the case for many state-of-the-art systems, this initial step tends to over-generate distinct named entities (several different entities are generated for the same real entity), favoring precision over recall. With this extra assumption, the cross-document co-reference task amounts to merge different entities into a single one, a problem that can be formulated as “alias detection”.

The event recognition and tracking task is solved by adapting incremental clustering techniques to deal with temporal constraints, implementing the intuitive idea that an event should emerge as a cluster of items reporting about a restricted set of actors (people, organization) interacting in a specific location during a certain time period. The in-

tuition behind our approach is that better cross-document named entity co-reference leads to more accurate similarity measures for clustering and better event recognition leads to better cross-document named entity co-reference. Thereby, we use the event profile of an entity as a powerful, synthetic representation of the socio-geo-temporal information related to the entity.

Our general algorithm is outlined in Alg. 1, and consists of successive applications of two steps: *clustering* and *co-reference*.

The clustering algorithm takes as input a set $\mathcal{D} = d_1 \dots d_n$ of documents, threshold parameter τ and a co-reference mapping S . The output P is a $n \times k$ matrix of cluster assignment. Typically, a variant of the fully-incremental (or single-pass) algorithm [17], described in Sect. 2, is used. Although in our experiments we performed an hard assignment, our proposed method can easily be extended to a soft-clustering setting that outputs as $P[i, j]$ the probability $p(c_j | d_i)$ that document d_i belongs to cluster c_j . As said above, the number of clusters k is not fixed but indirectly controlled by the threshold parameter τ that determines the granularity level of the final clustering. Each document d_i will be represented by three features: its bag-of-words, its bag-of-entities and its timestamp. The clustering algorithm then uses a similarity measures that combines these three features for each one of the documents. In our experiments, we weighted words and entities independently with a standard ‘tf-idf’ scheme, concatenated them into one vector and applied cosine similarity. Because an event is by nature limited in time, we placed an additional time constraint when comparing a document d to a cluster. Any cluster whose mean time (the average of the timestamps of the news articles composing it) exceeds 12 days from the timestamp of d is filtered out, i.e. considered as inactive. As pointed out by [18] this is done not only for computational efficiency, but also because similar articles separated in time are less likely to belong to the same event.

The co-reference mapping S decides which entities will be considered as co-referring. Its formal definition will be given later. It is used in the clustering algorithm as an inner entity-to-entity (or feature-to-feature) similarity matrix (like in the Generalized Vector Space model). More specifically, the S matrix is a matrix representing equivalence relations between features (entities) as discovered by the cross document named entity co-reference module.

Let us now turn to the cross document named entity co-reference task. The *coreference* algorithm takes as input the set of $\mathcal{E} = \{e_1, \dots, e_m\}$ entities and a clustering P of the document set \mathcal{D} , in addition to two parameters α and θ . In our experiments, we take a pairwise approach where we compute a similarity between entities and decide to merge entities e_i and e_j if $\text{sim}(e_i, e_j) \geq \theta$. This similarity is a convex linear combination between a *content* and a *context* similarity.

The *content* similarity only considers the canonical strings of the entities. Each entity e is represented through a sparse bag-of-words vector $w(e)$ with a non-zero entry for word v if the string representing entity e contains v . We weight this

again by a ‘tf-idf’ scheme, where the *idf* part here is the inverse of the number of distinct entities in which the word appears. As it is desirable to take into account some forms of fuzzy matching between the constitutive words (misspellings, transliterations, *etc.*), we also introduce a word-to-word similarity matrix Y , so that the content similarity measure between two entities e_i and e_j will be:

$$\text{sim}_{\text{content}}(e_i, e_j) = w(e_i)^T Y w(e_j) \quad (1)$$

where $w(e_i)$ and $w(e_j)$ are $L2$ -normalized. The matrix Y is derived here from the weighted-edit distance between two words. We learnt the weights of the edit distance using an external resource, namely the *JRC-Names* list of equivalent names, after alignment of their constitutive words³.

While the content similarity is fixed for two given entities, their *context* similarity varies depending on the clustering output. For each entity e we takes its event-profile, namely the binary vector $c(e)$ of size k whose i^{th} component is one if e appears in at least one document belonging the cluster i . Therefore,

$$\text{sim}_{\text{context}}(e_i, e_j) = c(e_i)^T c(e_j) \quad (2)$$

with $c(e_i)$ and $c(e_j)$ being $L2$ -normalized.

The final similarity is then a linear combination (late fusion) of these two measures:

$$\text{sim}(e_i, e_j) = \alpha \text{sim}_{\text{content}}(e_i, e_j) + (1 - \alpha) \text{sim}_{\text{context}}(e_i, e_j) \quad (3)$$

Recall that, in our case, the cross document named entity co-reference amounts to merge some entities together, as the intra document co-reference solver tends to be too specific. We therefore compute $\text{sim}(e_i, e_j)$ for all pairs of entities and merge those whose similarity is above θ . We then take the transitive closure of this relationship, obtaining the $m \times m$ symmetric matrix S with 1 at position $S[i, j]$ if e_i is merged with e_j . S is then used as parameter for the *clustering* module, controlling which of the entities are being considered as being equal.

Putting all this together, we formalize in Algorithm 1 our joint event recognition - cross-document co-reference algorithm, as the iterated application of two successive grouping algorithms, each one being influenced by the outcome of the previous one.

Algorithm 1 Joint Entity and Event Clustering

Input: a document set \mathcal{D} , the entity set \mathcal{E} , parameters τ, θ and α

- 1: $S^{(0)} = I$, the identity matrix
 - 2: $t = 1$
 - 3: **while** convergence criterion is not satisfied, **do**
 - 4: $P^{(t)} = \text{clustering}(\mathcal{D}, S^{(t-1)}, \tau)$
 - 5: $S^{(t)} = \text{coreference}(\mathcal{E}, P^{(t)}, \alpha, \theta)$
 - 6: $t = t + 1$
 - 7: **end while**
 - 8: **return** $P^{(t)}$ and $S^{(t)}$
-

³see <http://langtech.jrc.it/JRC-Names.html>

4. EXPERIMENTS AND DISCUSSION

For the experiments, we used the TDT5 news article corpus⁴. This collection includes about 280,000 documents (we considered only the articles written in English). Each document was linguistically pre-processed by the Xerox Incremental Parser [1] and this pre-processing included intra-document co-reference resolution. We limited ourselves to the following list of named entities: persons, places, organizations and linguistic events⁵. We represent each document by the concatenation of the tf-idf vectors of its named entities and bag-of-words (early fusion) or of its named entities alone (see below).

We evaluated the clustering using the TDT5 ground truth. In this dataset, 6364 articles are annotated with 126 events (called stories or topics in TDT5), which we took as gold reference for assessing the clustering performance. The clustering algorithm is the fully-incremental one, with time constraints (“old” clusters are filtered out) and including the entity co-reference mapping as explained in Section 3. Note that, while the algorithms are applied on the 280,000 documents, for assessing the clustering performance we extracted the subset corresponding to the 6364 labeled documents. We report micro-average precision and recall; adopting, as usual, the mapping between identified clusters and reference events that maximized the F_1 measure.

We fixed $\tau = 0.1$, a value that leads to the best results for the clustering at the first iteration (which will be used as a baseline).

In Figure 2 we plotted the evolution of the micro-average F_1 measure of the clustering over the different iterations. We ran this experiment with different values of θ and used two baselines:

1. The clustering after the first iteration
2. The clustering using words only (without considering named entities).

For the co-reference module, we fixed α manually to 0.75, a value which has not been optimized further. In order to see the impact of the social similarity, we run the experiments with $\alpha = 1$. In general, the results were slightly better than those obtained with $\alpha = 0.75$ at the first round, but worse than the results after convergence. With a value of $\alpha = 1$, it does not make sense to use a bootstrapping cycle like the one depicted in Figure 1, as the co-reference decision is independent of the current clustering.

Considering Figure 2, a first observation to note is the fact that the baseline using only the bag-of-words as features performs (slightly) better than the case where this is combined with the bag-of-entities.

However, after a couple of iterations improving the co-reference resolution, the clustering results also improve, emphasizing

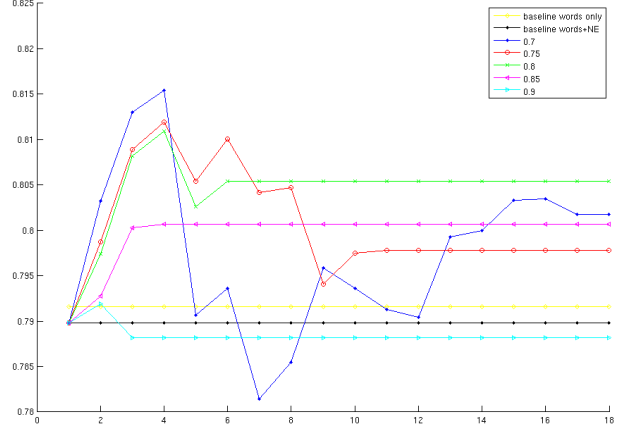


Figure 2: Evolution of the clustering (F_1) over successive iterations using words+NE

the importance of having a good entity resolution. As can be seen in this figure, taking a value of θ (the threshold for the co-reference module) between 0.75 and 0.85 improves the clustering in a consistent way. Outside of this range, the evolution is unstable, or even negative. Indeed, with a lower value of θ , the co-reference component precision decreases, introducing errors that affect negatively the next clustering iterations. With a too high value of θ , fewer referents are marked as being coreferring (reducing recall) and the changes are not significant enough to affect the clustering results.

In order to further investigate the connections between having a better co-reference resolution of entities and the clustering quality, we represented each document as the vector of its named entities alone, instead of combining the words and the named entities. The corresponding results on the clustering are plotted in Figure 3. As expected, event detection performance decreases as the entities alone do not seem to carry enough information to cluster correctly the documents. While some of the previous conclusions also hold in this case (a general improvement), the results seem to be much more sensitive to the chosen value of θ . This seems to indicate that the output of the co-reference module is not very stable, which can be explained by seeing the low values for the clustering ($F_1 = 0.56$ is a rather low value for a clustering). The fact that a relatively bad input produces inconsistent results should not come as a surprise.

Until now, we only evaluated the co-reference of named entities indirectly, through its impact on the clustering. In order to evaluate the opposite (the impact of the clustering on the co-reference resolution) directly, we created our own gold standard. From the 126 events provided by the TDT5 ground truth, we took randomly 56 of these events and analyzed the named entities that were detected inside articles of these events. There were in total 4516 different named entities. We manually annotated all the reference pairs that refer to the same real entity. Note that, in order to reduce the amount of comparison, we considered only pairs that ap-

⁴www.ldc.upenn.edu/Projects/TDT5

⁵contrary to what we call “event” in this paper, a linguistic event is a lexical expression referring to an event such as “London Olympic Games”

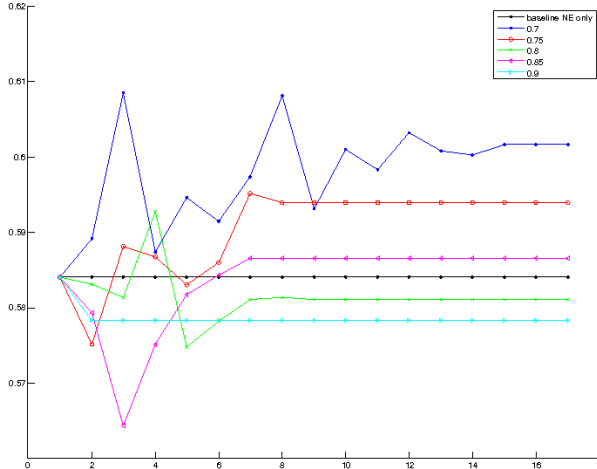


Figure 3: Evolution of the clustering (F_1) over successive iterations using NE only

peared in the same event. Those annotated pairs that were referring to the same underlying entity constituted our cross-document named entity co-reference ground truth. This resulted in 591 equivalence classes that contained more than one named entity. The fact of considering only entities that co-occur in at least one of the selected events is less than optimal: we do not consider this to be a clean gold standard, but rather a sanity check to see if the performance of both tasks are effectively correlated.

In Figure 4 we plotted the evolution of the micro-average F_1 measure of successive results of the co-reference module, with respect to our annotated ground truth, using $\theta = 0.8$ during the process, the threshold that obtained the best result for the clustering. Note that the values plotted there correspond actually to the threshold that achieved best F_1 (this was always around 0.55): the reason why we used a larger threshold in the actual process was to introduce fewer false positives, in exchange of introducing fewer changes in general. As can be appreciated, the shape of the figure mirrors the corresponding curve in Figure 2. This illustrates that for each iteration, the clustering is strongly linked to the co-reference resolution result, adding evidence to our initial assumption that entity resolution and event detection are effectively correlated. At the same time, the general low performance achieved, shows that the idea can still be improved. In particular, considering the cross-document co-reference as an “alias detection task” works for the first iterations but after the fourth iteration it seems to create too many aliasing with the consequence of decreasing the result quality.

5. CONCLUSIONS

We considered two apparently independent problems in a larger application (a news aggregator) and tackled them jointly. We adapted the algorithms used for the clustering of events and for co-referencing resolution of named entities to allow as additional input the output of the other algorithm.

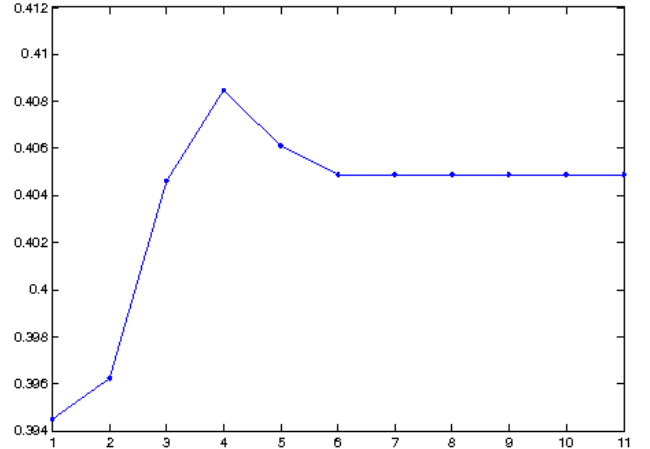


Figure 4: Evolution of the cross-document co-reference (F_1) over successive iterations for $\theta = 0.8$.

This feedback procedure can then be used in a bootstrap fashion to improve both tasks simultaneously. Such a general bootstrapping approach is of general interest and may very probably be applied to a wide range of applications for data mining on the web.

In the experiments performed for this paper we limited the number of features for the co-reference resolution. The evaluation results for the co-reference module (remind that our ground truth was more of a guideline) were rather low, but still had a positive impact on the resulting clustering. State-of-the-art co-reference systems use up to 60 different features (see Sect. 2): adding some of those and weighting them properly should help the co-reference resolution performance, and by extension, the clustering. The point of this paper is not to present a state-of-the-art co-reference module, but to show the improvement that adding social features can have in such a task. In order to reduce the number of parameters we only used one content feature, but the literature in this field confirms that choosing linguistically richer features could potentially improve the co-reference. In particular we did not add additional features for acronyms, which are recurrent in the news-sphere. Of course, this would need a proper (supervised) learning algorithm to learn the parameters, as for now we fixed the only parameter (α) by hand.

Related to this, our general approach for co-referencing entities is based on pairwise similarities. Previous studies have shown that directly considering a chain structure [9] or other richer structures [8, 10] improves the result. One advantage of our proposed method is that any feature-based co-reference module can be used, including those we just mentioned. By just adding the context similarity as an additional feature to the mentioned methods, previous work could be plugged into our framework.

Our cross-document co-reference module corrects errors where two expressions point to the same entity but does not consider cases of ambiguous entities. Moreover, we assumed

that the canonical forms returned by the initial co-reference detection engine are as specific as possible. While our experience showed us that this is a valid assumption in most cases, there are cases where the expressions are not specific enough. Obvious cases in the news sphere (where new entities are continually appearing) are of course people who share exactly the same name, but other cases can happen also. Take as example an article where all references of “New York” are actually referring to the city of New York. Merging this with the expression “New York City” would be correct in this context, but wrong in a context where this same string refers to the US state instead of the city. We believe that this could explain a strange phenomenon that appears in Figure 2 and 4: a sharp jump in the first iterations, followed by a drop and a slower ascent. This may be due to the fact that at the first iterations we do correct most of the errors due to overspecify, but then perform more errors by merging entities that do not point to the same real entity. An improvement to this would be a system where the merging of expressions would be done in a context-specific way, instead of globally. Each event may have a co-reference table that indicates which expression lead to the same entity *in this context* and new documents will be compared using that table. This does however underuses the available information by restricting only to the context where enough evidence is found. It would be desirable to define a notion of similar contexts in order to generalize the obtained information. Related to this is the consideration of a flexible threshold θ that changes at each iteration.

6. REFERENCES

- [1] S. Ait-Mokhtar, J.-P. Chanod, and C. Roux. Robustness beyond shallowness: incremental deep parsing. *Nat. Lang. Eng.*, 8:121–144, June 2002.
- [2] J. Allan, R. Papka, and V. Lavrenko. On-line new event detection and tracking. In *SIGIR*, pages 37–45, New York, New York, USA, 1998. ACM Press.
- [3] I. Bhattacharya and L. Getoor. Relational clustering for multi-type entity resolution. In *Proceedings of the 4th international workshop on Multi-relational mining*, MRDM ’05, pages 3–12, New York, NY, USA, 2005. ACM.
- [4] I. Bhattacharya and L. Getoor. Collective entity resolution in relational data. *ACM Trans. Knowl. Discov. Data*, 1(1), Mar. 2007.
- [5] M. Gallé and J.-M. Renders. Full and semi-batch clustering of news articles with Star-EM. In *34th European Conference on IR Research, ECIR 2012*. Springer, 2012.
- [6] J. Huang, P. Treeratpituk, S. M. Taylor, and C. L. Giles. Enhancing cross document coreference of web documents with context similarity and very large scale text categorization. In *Proceedings of the 23rd International Conference on Computational Linguistics*, COLING ’10, pages 483–491, Stroudsburg, PA, USA, 2010. Association for Computational Linguistics.
- [7] A. Kotov, C. Zhai, and R. Sproat. Mining named entities with temporally correlated bursts from multilingual web news streams. In *Proceedings of the Forth International Conference on Web Search and Web Data Mining, WSDM 2011, Hong Kong, China, February 9-12, 2011*, pages 237–246. ACM, 2011.
- [8] X. Luo, A. Ittycheriah, H. Jing, N. Kambhatla, and S. Roukos. A mention-synchronous coreference resolution algorithm based on the bell tree. In *Proceedings of the Annual Meeting on Association for Computational Linguistics*, page 135. Association for Computational Linguistics, 2004.
- [9] V. Ng and C. Cardie. Improving machine learning approaches to coreference resolution. In *Proceedings of the Annual Meeting on Association for Computational Linguistics*, page 104, Morristown, NJ, USA, July 2002. Association for Computational Linguistics.
- [10] C. Nicolae and G. Nicolae. BestCut: A Graph Algorithm for Coreference Resolution. In *Proceedings of Empirical Methods in Natural Language Processing*, pages 275–283, July 2006.
- [11] T. Pedersen, A. Purandare, and A. Kulkarni. Name discrimination by clustering similar contexts. In *Proceedings of Conference on Intelligent Text Processing and Computational Linguistics*, pages 226–237. Springer, 2004.
- [12] B. Poulighen, R. Steinberger, C. Ignat, I. Temnikova, A. Widiger, W. Zaghouni, and J. Zizka. Multilingual person name recognition and transliteration. *CORELA - Cognition, Representation, Language*, 3(2):115–123, 2006.
- [13] Y. Shinyama and S. Sekine. Named entity discovery using comparable news articles. In *Proceedings of International Conference on Computational Linguistics*, Stroudsburg, PA, USA, 2004. Association for Computational Linguistics.
- [14] R. Steinberger, B. Poulighen, M. Kabadjov, J. Belyaeva, and E. van der Goot. JRC-NAMES: A Freely Available, Highly Multilingual Named Entity Resource. In *Proceedings of Recent Advances in Natural Language Processing*, pages 14–16, 2011.
- [15] V. Stoyanov, N. Gilbert, C. Cardie, and E. Riloff. Conundrums in Noun Phrase Coreference Resolution : Making Sense of the State-of-the-Art. In *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP: Volume 2-Volume 2*, pages 656–664. Association for Computational Linguistics, 2009.
- [16] S. Vadrevu, C. Teo, S. Rajan, K. Punera, B. Dom, A. J. Smola, Y. Chang, and Z. Zheng. Scalable clustering of news search results. In *WSDM*, pages 675–683, 2011.
- [17] C. van Rijsbergen. *Information Retrieval*. Butterworths, 1979.
- [18] Y. Yang, T. Pierce, and J. Carbonell. A study of retrospective and on-line event detection. In *SIGIR*, pages 28–36. ACM, 1998.